

# Случайные процессы. Прикладной поток.

## Практическое задание 3

### Правила:

- Выполненную работу нужно отправить на почту `probability.diht@yandex.ru`, указав тему письма "[СП17] Фамилия Имя - Задание 3". Квадратные скобки обязательны. Вместо Фамилия Имя нужно подставить свои фамилию и имя.
- Прислать нужно ноутбук и его pdf-версию. Названия файлов должны быть такими: `3.N.ipynb` и `3.N.pdf`, где `N` - ваш номер из таблицы с оценками.
- При проверке некоторый код из вашего решения будет **проверяться автоматически**. Этот код вы должны скопировать в файл с названием `s3.N.py` и прислать вместе с решением. Что именно должно быть в этом файле, написано далее. Код должен корректно работать в Python 3.5 под Убунту.

---

Для выполнения задания потребуются следующие библиотеки: `bs4`, `urllib`, `networkx`. Следующими командами можно их поставить (Ubuntu):

```
sudo pip3 install beautifulsoup4
```

```
sudo pip3 install urllib2
```

```
sudo pip3 install networkx
```

---

## PageRank

### История

(Взято с Википедии (<https://ru.wikipedia.org/wiki/PageRank>))

В 1996 году Сергей Брин и Ларри Пейдж, тогда ещё аспиранты Стэнфордского университета, начали работу над исследовательским проектом BackRub — поисковой системой по Интернету, использующей новую тогда идею о том, что веб-страница должна считаться тем «важнее», чем больше на неё ссылаются других страниц, и чем более «важными», в свою очередь, являются эти страницы. Через некоторое время BackRub была переименована в Google. Первая статья с описанием применяющегося в ней метода ранжирования, названного PageRank, появилась в начале 1998 года, за ней следом вышла и статья с описанием архитектуры самой поисковой системы.

Их система значительно превосходила все существовавшие тогда поисковые системы, и Брин с Пейджем, осознав её потенциал, основали в сентябре 1998 года компанию Google Inc., для дальнейшего её развития как коммерческого продукта.

### Описание

Введем понятие веб-графа. Ориентированный граф  $G = (V, E)$  называется веб-графом, если

- $V = \{url_i\}_{i=1}^n$  --- некоторое подмножество страниц в интернете, каждой из которых соответствует адрес  $url_i$
- Множество  $E$  состоит из тех и только тех пар  $(url_i, url_j)$ , для которых на странице с адресом  $url_i$  есть ссылка на  $url_j$

Рассмотрим следующую модель поведения пользователя. В начальный момент времени он выбирает некоторую страницу из  $V$  в соответствии с некоторым распределением  $\Pi^{(0)}$ . Затем, находясь на некоторой странице, он может либо перейти по какой-то ссылке, которая размещена на этой странице, либо выбрать случайную страницу из  $V$  и перейти на нее (damping factor). Считается, что если пользователь выбирает переход по ссылке, то он выбирает равновероятно любую ссылку с данной страницы и переходит по ней. Если же он выбирает переход не по ссылке, то он также выбирает равновероятно любую страницу из  $V$  и переходит на нее (в частности может остаться на той же странице). Будем считать, что переход не по ссылке пользователь выбирает с некоторой вероятностью  $p \in (0, 1)$ . Соответственно, переход по ссылке он выбирает с вероятностью  $1 - p$ . Если же со страницы нет ни одной ссылки, то будем считать, что пользователь всегда выбирает переход не по ссылке.

Описанная выше модель поведения пользователя называется моделью PageRank. Нетрудно понять, что этой модели соответствует некоторая марковская цепь. Опишите ее.

- Множество состояний: <Ответ>
- Начальное распределение: <Ответ>
- Переходные вероятности: <Ответ (формула)>

## Вычисление

Данная марковская цепь является эргодической. Почему?

<Ответ>

А это означает, что цепь имеет некоторое эргодическое распределение  $\Pi$ , которое является предельным и единственным стационарным. Данное распределение называется весом PageRank для нашего подмножества интернета.

Как вычислить это распределение  $\Pi$  для данного веб-графа? Обычно для этого используют степенной метод (power iteration), суть которого состоит в следующем. Выбирается некоторое начальное распределение  $\Pi^{(0)}$ . Далее производится несколько итераций по формуле  $\Pi^{(k)} = \Pi^{(k-1)} P$ , где  $P$  --- матрица переходных вероятностей цепи, до тех пор, пока  $\|\Pi^{(k)} - \Pi^{(k-1)}\| > \varepsilon$ . Распределение  $\Pi^{(k)}$  считается приближением распределения  $\Pi$ .

Имеет ли смысл выполнять подобные итерации для разных начальных распределений  $\Pi^{(0)}$  с точки зрения теории?

<Ответ>

А с точки зрения практического применения, не обязательно при этом доводя до сходимости?

<Ответ>

Какая верхняя оценка на скорость сходимости?

<Ответ>

# Часть 1

In [ ]:

```
import numpy as np
from scipy.stats import bernoulli
import networkx
from bs4 import BeautifulSoup
from urllib.request import urlopen
from urllib.parse import urlparse, urlunparse
from time import sleep
from itertools import product
import matplotlib.pyplot as plt

%matplotlib inline
```

Реализуйте вычисление весов PageRank power-методом.

Реализовать может быть удобнее с помощью функции `np.nan_to_num`, которая в данном `numpy.array` заменит все вхождения `nan` на ноль. Это позволяет удобно производить поэлементное деление одного вектора на другой в случае, если во втором векторе есть нули.

In [ ]:

```
def create_page_rank_markov_chain(links, damping_factor=0.15):
    ''' По веб-графу со списком ребер links строит матрицу
    переходных вероятностей соответствующей марковской цепи.

    links --- список (list) пар вершин (tuple),
               может быть передан в виде numpy.array, shape=(|E|, 2);
    damping_factor --- вероятность перехода не по ссылке (float);

    Возвращает prob_matrix --- numpy.matrix, shape=(|V|, |V|).
    ...

    links = np.array(links)
    N = links.max() + 1 # Число веб-страниц

    <...>

    return prob_matrix

def page_rank(links, start_distribution, damping_factor=0.15,
              tolerance=10 ** (-7), return_trace=False):
    ''' Вычисляет веса PageRank для веб-графа со списком ребер links
    степенным методом, начиная с начального распределения start_distribution,
    доводя до сходимости с точностью tolerance.

    links --- список (list) пар вершин (tuple),
               может быть передан в виде numpy.array, shape=(|E|, 2);
    start_distribution --- вектор размерности |V| в формате numpy.array;
    damping_factor --- вероятность перехода не по ссылке (float);
    tolerance --- точность вычисления предельного распределения;
    return_trace --- если указана, то возвращает список распределений во
                     все моменты времени до сходимости

    Возвращает:
    1). если return_trace == False, то возвращает distribution ---
    приближение предельного распределения цепи,
    которое соответствует весам PageRank.
    Имеет тип numpy.array размерности |V|.
    2). если return_trace == True, то возвращает также trace ---
    список распределений во все моменты времени до сходимости.
    Имеет тип numpy.array размерности
    (количество итераций) на |V|.
    ...

    prob_matrix = create_page_rank_markov_chain(links,
                                                damping_factor=damping_factor)
    distribution = np.matrix(start_distribution)

    <...>

    if return_trace:
        return np.array(distribution).ravel(), np.array(trace)
    else:
        return np.array(distribution).ravel()
```

Автоматическая проверка

Реализацию функций `create_page_rank_markov_chain` и `page_rank` скопируйте в файл с названием `с3.N.py` и вышлите на почту. Будет проверяться только корректность выдаваемых значений. Проверки на время работы не будет.

---

Давайте посмотрим, как оно работает. Напишите для начала функцию для генерации случайного ориентированного графа  $G(n, p)$ . Случайный граф генерируется следующий образом. Берется множество  $\{0, \dots, n - 1\}$ , которое есть множество вершин этого графа. Ребро  $(i, j)$  (пара упорядочена, возможно повторение) добавляется в граф независимо от других ребер с вероятностью  $p$ .

In [ ]:

```
def random_graph(n, p):  
    return <Список ребер. Сможете в одну строчку? ;)>
```

Теперь сгенерируем случайный граф и нарисуем его.

In [ ]:

```
N, p = 10, 0.2  
edges = random_graph(N, p)  
  
G = networkx.DiGraph()  
G.add_edges_from(edges)  
plt.axis('off')  
networkx.draw_networkx(G, width=0.5)
```

Посчитаем его PageRank и изобразим так, чтобы размер вершины был пропорционален ее весу.

In [ ]:

```
start_distribution = np.ones((1, N)) / N  
pr_distribution = page_rank(edges, start_distribution)  
  
size_const = 10 ** 4  
plt.axis('off')  
networkx.draw_networkx(G, width=0.5, node_size=size_const * pr_distribution,  
                        node_color=pr_distribution)
```

Как мы уже отмечали выше, эргодическая теорема дает верхнюю оценку на скорость сходимости. Давайте посмотрим, насколько она является точной. Для этого при вычислении PageRank нужно установить флаг `return_trace`.

In [ ]:

```
pr_distribution, pr_trace = page_rank(edges, start_distribution,
                                     return_trace=True)
errors = np.abs(pr_trace - pr_trace[-1]).sum(axis=(1, 2))

plt.figure(figsize=(10, 4))
x = np.arange(len(errors))
plt.plot(x, errors, lw=2, label='error')
plt.plot(x, <верхняя оценка скорости сходимости из эргодической теоремы>,
         lw=2, label='estimation')
plt.legend()
plt.xlabel('iterations')
plt.show()
```

<Выводы>

Проведите небольшое исследование. В ходе исследования выясните, как скорость сходимости (количество итераций до сходимости) зависит от  $n$  и  $p$ , а так же начального распределения. Вычислите также веса PageRank для некоторых неслучайных графов. В каждом случае стройте графики. От чего зависит вес вершины?

<Исследования и выводы>

## Часть 2

В этой части вам предстоит построить реальный веб-граф и посчитать его PageRank. Ниже определены вспомогательные функции.

In [ ]:

```
def load_links(url, sleep_time=1, attempts=5, timeout=20):
    ''' Загружает страницу по ссылке url и выдает список ссылок,
    на которые ссылается данная страница.
        url --- string, адрес страницы в интернете;
        sleep_time --- задержка перед загрузкой страницы;
        timeout --- время ожидания загрузки страницы;
        attempts --- число попыток загрузки страницы.
        Попытка считается неудачной, если выбрасывается исключение.

        В случае, если за attempts попыток не удалось загрузить страницу,
        то последнее исключение пробрасывается дальше.
    '''

    sleep(sleep_time)
    parsed_url = urlparse(url)
    links = []

    # Попытки загрузить страницу
    for i in range(attempts):
        try:
            # Ловить исключения только из urlopen может быть недостаточно.
            # Он может выдавать какой-то бред вместо исключения,
            # из-за которого исключение сгенерирует BeautifulSoup
            soup = BeautifulSoup(urlopen(url, timeout=timeout), 'lxml')
            break

        except Exception as e:
            print(e)
            if i == attempts - 1:
                raise e

    for tag_a in soup('a'): # Посмотр всех ссылочных тегов
        if 'href' in tag_a.attrs:
            link = list(urlparse(tag_a['href']))

            # Если ссылка является относительной,
            # то ее нужно перевести в абсолютную
            if link[0] == '': link[0] = parsed_url.scheme
            if link[1] == '': link[1] = parsed_url.netloc

            links.append(urlunparse(link))

    return links

def get_site(url):
    ''' По ссылке url возвращает адрес сайта. '''

    return urlparse(url).netloc
```

Код ниже загружает  $M$  веб-страниц, начиная с некоторой стартовой страницы и переходя по ссылкам. Загрузка происходит методом обхода в ширину. Все собранные урлы страниц хранятся в `urls`. В `links` хранится список ссылок с одной страницы на другую. Особенность кода такова, что в `urls` хранятся все встреченные урлы, которых может быть сильно больше  $M$ . Аналогично, в `links` ребра могут ссылаться на страницы с номером больше  $M$ . Однако, все ребра из `links` начинаются только в первых  $M$  страницах. Таким образом, для построения веб-графа нужно удалить все, что связано с вершинами, которые не входят в первые  $M$ .

Это очень примерный шаблон, к тому же не оптимальный. Можете вообще его не использовать и написать свое.

In [ ]:

```
urls = ['http://wikipedia.org/wiki/']
site = get_site(urls[0])
links = []

N = 10
for i in range(N):
    try:
        # Загружаем страницу по урлу и извлекаем из него все ссылки
        # Не выставляйте sleep_time слишком маленьким,
        # а то еще забанят где-нибудь
        links_from_url = load_links(urls[i], sleep_time=0.5)
        # Если мы хотим переходить по ссылкам только определенного сайта
        links_from_url = list(filter(lambda x: get_site(x) == site,
                                     links_from_url))

        # Добавляем соответствующие вершины и ребра в веб-граф
        for j in range(len(links_from_url)):
            # Такая ссылка уже есть
            if links_from_url[j] in urls:
                links.append((i, urls.index(links_from_url[j])))

            # Новая ссылка
            else:
                links.append((i, len(urls)))
                urls.append(links_from_url[j])

    except:
        pass # Не загрузилась с 5 попытки, ну и ладно
```

Теперь выберите какой-нибудь сайт с небольшим количеством страниц (не более 1000). Таким сайтом может быть, например, сайт кафедры Дискретной математики (<http://ru.discrete-mathematics.org>) (аккуратнее, если забанят, то лишитесь доступа к учебным материалам ;), Школы анализа данных (<http://yandexdataschool.ru>), сайт магазина, больницы.

Постройте полный веб-граф для этого сайта и визуализируйте его. При отрисовке выставляйте width не более 0.1, иначе получится ужасно некрасиво.

Посчитайте PageRank для этого веб-графа. Визуализируйте данный веб-граф, сделав размер вершин пропорционально весу PageRank (см. пример в части 1). Постройте гистограмму весов. Что можно сказать про скорость сходимости?

Выделите небольшое количество (15-20) страниц с наибольшим весом и изобразите граф, индуцированный на этом множестве вершин. Что это за страницы? Почему именно они имеют большой вес?

Как меняется вес PageRank для страниц в зависимости от начального приближения в случае, если не доводить итерационный процесс вычисления до сходимости? Какие выводы о поведении пользователя отсюда можно сделать?

Для получения дополнительных баллов проведите аналогичные исследования для больших сайтов. Также вы можете провести исследования, не ограничиваясь загрузкой только одного сайта.